

METHOD FOR GENERATING FIVE PRIME BIASED TANDEM TAG LIBRARIES OF cDNAs

BACKGROUND OF THE INVENTION

1. FIELD OF THE INVENTION

The sequences of whole genomes from several organisms have now been elucidated and are available as searchable databases. This enables rapid identification of full-length messenger RNAs (mRNAs) expressed in a biological sample once a partial sequence is known. The method described here allows generation of such partial sequences consisting of a minimal length of expressed cDNA sequences of at least 20 bases from biological samples to rapidly identify novel expressed transcripts.

2. DESCRIPTION OF THE RELATED ART

In order to obtain a comprehensive collection of all human genes that are expressed, many millions of cDNA molecules must be sequenced, which is quite costly and laborious. Since the availability of the human genome sequence, much of the coding sequence of a gene can now be inferred once a short physical sequence is obtained. Hence, sequencing only a short stretch of cDNAs should be sufficient in theory to identify all genes expressed in a biological sample. The Expressed Sequence Tag (EST) method

10092286 "030602

U.S. PATENT APPLICATION OF SAMAL ET AL.

purports to achieve this by generating for sequencing relatively short cDNA fragments from 3' ends. However, the EST method still utilizes one cDNA per clone, which means one sequencing reaction yields one cDNA sequence.

An effective way to improve this yield so that each plasmid and each sequencing reaction yields many cDNA sequences is to "glue" together short cDNA fragments from end to end. The Serial Gene Expression Analysis (SAGE) method effectively utilizes such a concatenation procedure. The SAGE method, however, has two key shortcomings. One is that all of the tags are generated from a defined 3' end of a cDNA. Mammalian genes contain long untranslated sequences at their 3' ends, which make the determination of coding sequence by gene prediction algorithms difficult and unreliable. The second limitation is that the SAGE tags are typically only 14 bases long, which are too short to yield uniquely matching sequences from the genomic database. A minimum of 20 bases is needed to identify a uniquely matching gene from a mammalian genomic database at 80% of the time.

Thus, the most important prerequisite for obtaining expressed sequence tags to rapidly and uniquely identify coding sequences from a messenger RNA pool is to obtain

U.S. PATENT APPLICATION OF SAMAL ET AL.

expressed sequence tags of 20 bases or longer from the 5' end of a coding region. Such tags then can be used as a forward PCR primer to easily amplify, sequence, and clone each gene uniquely. There is presently no method, which predictably generates 5' cDNA fragments of 20-40 bases. The method described here generates one or more short tags at or near the 5' end of each gene transcript in tandem or in cluster so that when they are aligned against genomic sequences they together uniquely identify a contiguous expressed sequence of 20 bases or greater.

SUMMARY OF THE INVENTION

The present application discloses a method for generating five prime biased tandem tag libraries of cDNAs. The method comprises the steps of isolating a sample of mRNAs; synthesizing double-stranded cDNAs from the mRNAs; blunt-ending the double-stranded cDNAs; attaching an adapter molecule to the blunt ends of the double stranded cDNAs to form a complex, where the adapter molecule is a double stranded, synthetic oligonucleotide comprising a recognition site for a type IIS restriction enzyme, a cloning site for releasing tags to a cloning vector, and

[illegible]

In a preferred embodiment, the type IIS restriction enzyme is selected from the group consisting of Ear I, Sap I, Alw I, Bmr I, Bsa I, BsmA I, BsmB I, Mly I, Ple I, Bbs I, BciV I, Fau I, Mnl I, Aar I, BfuA I, BspM I, Hph I, Mbo II, SspD5 I, Sth132 I, SfaN I, BseR I, BspCN I, Hga I, AceIII, Eci I, TaqII, Tth111III, Bbv I, RleAI, Bcefi, Fok I, BceA I, BsmF I, StsI, Bce83I, BpmI, Bsg I, Eco57I, Eco57MI, and MmeI. In a more preferred embodiment, the type IIS restriction enzyme is BpmI.

In other preferred embodiments, the released tags are comprised of 50 nucleotides or less; the released tags are comprised of 36 nucleotides or less; the released tags are comprised of 32 nucleotides or less. In a more preferred

U.S. PATENT APPLICATION OF SAMAL ET AL.

embodiment, the released tags are comprised of at least 20 nucleotides.

In yet another preferred embodiment, the method further comprises sequencing the isolated concatenated tags to obtain a nucleotide sequence and comparing the nucleotide sequence to a known nucleotide sequence.

The present application also discloses a method for generating five prime biased tandem tag libraries of cDNAs, comprising the steps of isolating a sample of mRNAs; synthesizing double-stranded cDNAs from the mRNAs; blunt-ending the double-stranded cDNAs; attaching a first adapter molecule to the blunt ends of the double stranded cDNAs to form a first complex, where the first adapter molecule is a double stranded, synthetic oligonucleotide comprises a recognition site for a type IIS restriction enzyme, a cloning site for releasing tags to a cloning vector, and a PCR primer site; digesting the first complex with a type IIS restriction enzyme to form first released tags; separating the first released tags from the double-stranded cDNAs and attaching a second adapter molecule to the double-stranded cDNAs to form a second complex; amplifying the first released tags to form first amplified tags; isolating the first amplified tags; concatenating the first

U.S. PATENT APPLICATION OF SAMAL ET AL.

amplified tags to form first concatenated tags; amplifying the first concatenated tags; isolating the first concatenated tags; digesting the second complex with a type IIS restriction enzyme to form second released tags; separating the second released tags from the double-stranded cDNAs; amplifying the second released tags to form second amplified tags; isolating the second amplified tags; concatenating the second amplified tags to form second concatenated tags; amplifying the second concatenated tags; and isolating the second concatenated tags.

In a preferred embodiment, the type IIS restriction enzyme is selected from the group consisting of Ear I, Sap I, Alw I, Bmr I, Bsa I, BsmA I, BsmB I, Mly I, Ple I, Bbs I, BciV I, Fau I, Mnl I, Aar I, BfuA I, BspM I, Hph I, Mbo II, SspD5 I, Sth132 I, SfaN I, BseR I, BspCN I, Hga I, AceIII, Eci I, TaqII, Tth111II, Bbv I, RleAI, BceFI, Fok I, BceA I, BsmF I, StsI, Bce83I, BpmI, Bsg I, Eco57I, Eco57MI, and MmeI. In a more preferred embodiment, the type IIS restriction enzyme is BpmI.

In another preferred embodiment, the mRNAs are from a mammal. In a more preferred embodiment, the mRNAs are from a human.

U.S. PATENT APPLICATION OF SAMAL ET AL.

In other preferred embodiments, the released tags are comprised of 50 nucleotides or less; the released tags are comprised of 36 nucleotides or less; the released tags are comprised of 32 nucleotides or less. In a more preferred embodiment, the released tags are comprised of at least 20 nucleotides.

In yet another preferred embodiment, the method further comprises sequencing the isolated concatenated tags to obtain a nucleotide sequence and comparing the nucleotide sequence to a known nucleotide sequence.

BRIEF DESCRIPTION OF THE DRAWINGS

Figures 1A, 1B and 1C show a flow chart of an embodiment of the present method for generating five primed biased tandem tag libraries of cDNAs.

DETAILED DESCRIPTION

A. Brief Description of the Method

1. The first and second strand cDNA synthesis is carried out according the standard procedure. In a preferred embodiment, the first strand synthesis is carried out with oligo-dT 3' primer covalently linked to magnetic

U.S. PATENT APPLICATION OF SAMAL ET AL.

beads according to the manufacturer's protocol (Dyna
Inc.).

2. The 5' ends of the ds-cDNAs are flushed using T4 DNA polymerase in the presence of dNTP, followed by the ligation of a double stranded adaptor. The adaptor can be of any sequence but contains the recognition sequence for a type IIS restriction enzyme that cleaves double stranded DNA substrates at some length downstream of the recognition site. In a preferred embodiment, the recognition sequence for a type IIS enzyme, Bpm I (also known as Gsu I) was placed at the 3' end of the adaptor so that the nucleotide sequence immediately following the Bpm I site is from cDNAs. In addition, optionally, the recognition site for a rare six cutter such as the Mlu I enzyme can also be incorporated into the adaptor at just upstream of the Bpm I site to be utilized at a later step. The remaining adaptor sequence serves as the forward primer site for a subsequent PCR amplification step.

3. The ligated adaptor-cDNAs are purified and then digested with Bpm I to release the 16bp cDNA tags plus the adaptor. The rest of the cDNAs remain bound to the magnetic beads and saved.

U.S. PATENT APPLICATION OF SAMAL ET AL.

4. The adaptor-tag fragments are recovered by separating away the magnetic beads. They are ligated with a second adaptor of an arbitrary sequence but containing the same Mlu I site at the 5' end of the adaptor. These two adaptors also facilitate PCR amplification of the internal 16 bp cDNA tags.

5. PCR amplification is carried out according to the standard procedure using the forward and reverse primers, which contain the sequences of the two adaptors respectively. The product is purified and ligated to a PCR cloning vector followed up by the transformation of competent bacteria.

6. Plasmid harboring colonies are drug-selected. The plasmid DNA is purified and digested with Mlu I. The released tags plus the restriction sites (28 bp) are isolated and ligated to form concatamers. The concatmers of appropriate size, typically 0.5 Kb - 1.5 Kb, are fractionated by agarose gel-electrophoresis and then ligated into a Mlu I cut vector. After cloning, the 16 bp cDNA tags are elucidated by sequencing the concatemers.

7. The remaining cDNAs bound to the magnetic beads from the step 3 are then processed again through steps 2 - 6 to generate the second 16 bp tag from each cDNA. Thus,

U.S. PATENT APPLICATION OF SAMAL ET AL.

after the two rounds, two tandem tags from the 5' end of each expressed transcript are generated, which, when aligned against the genomic sequence, generate 32 bases of combined sequence.

8. Steps 2 - 6 can be repeated several times as necessary.

B. More Detailed Description of the Method

Step 1: cDNA synthesis

Total RNA was isolated from the HK 532 Cortical Cell Line using the Qiagen total RNA isolation kit (Qiagen, Inc., Valencia, CA). Briefly, the cells were lysed in a lysis buffer followed by binding of the RNA to the Qiagen solid matrix, from which the RNA was eluted, precipitated and kept at -20 °C overnight.

Messenger RNA (mRNA), typically of 200 ng, was incubated with Dynal beads (Dynal, Inc., Lake Success, NY) containing oligo(dT) to attach the polyadenylated RNA which was converted into cDNA using the Superscript II cDNA synthesis kit (GIBCO Life Technologies, Gaithersburg, MD) according to the manufacturer's directions.

Step 2: Adaptor ligation

After the second strand synthesis, the 5' ends of the double stranded-cDNA (ds-cDNA) were flushed using T4 DNA

U.S. PATENT APPLICATION OF SAMAL ET AL.

polymerase. Oligonucleotide adaptors were created by mixing equimolar amount of each of two synthetic oligonucleotides

sense strand:

GCAGTGGTATCAACGCAGAGTCCAGTGTGGTGGACGCGTCTGGAG (SEQ ID NO:1)

antisense strand:

pCTCCAGACGCGTCCACCACACTGGACTCTGCGTTGATACCAC (SEQ ID NO:2)

in deionized water, heating them to 95 °C, and allowing them to cool slowly to room temperature to form:

PCR primer site

MluI_BpmI

5' GCAGTGGTATCAACGCAGAGTCCAGTGTGGTGG**ACGCGT**CTGGAG
|||||
CACCATAGTTGCGTCTCAGGTCACACCACCT**GC**GCAGACCTC_p

(SEQ ID NO:3)

Adaptor DNA (500 pmoles) was added to the solid-phase cDNA in a total volume of 50 µl of 1× ligase buffer containing 25 U of T4 ligase (Gibco BRL). The reaction was incubated overnight at 16 °C followed by 10 min at 65 °C to inactivate the enzyme.

Step 3: Release and recovery of the first tag

Beads were again washed extensively in wash buffer (5 mM TrisHCl, pH 8.0, 0.5 mM EDTA, 1M NaCl and 200 µg

U.S. PATENT APPLICATION OF SAMAL ET AL.

BSA/ml), followed by three washes in BpmI buffer, and resuspended in 50 µl of Bpm I buffer containing 50 U of Bpm I and incubated at 37 °C for 5 h with gentle rotation. The tag-containing supernatant was collected and the beads were washed twice with 100 µl of reaction buffer 3 (NEBL, Beverly, MA). The supernatant and washes were combined. The combined material was extracted with phenol:CIA. A half volume of 7.5 M ammonium acetate, or a one-third volume of 10 M ammonium acetate was added and DNA was precipitated with 2 volumes of ethanol in the presence of 4 µl of glycogen (20 mg/ml) per 300 µl of initial volume.

Step 4: Ligation of the 3' adaptor

A second, 16-fold degenerate adaptor molecule was prepared by annealing synthetic oligos as described above

sense strand:

ACGCGTGTCTGACCTCGAGT (SEQ ID NO:4);

antisense strand:

TCTAGACTCGAGGTCGACACGCGTNN (SEQ ID NO:5)

to give the following oligodimer:

U.S. PATENT APPLICATION OF SAMAL ET AL.

Mlu I PCR primer site

^PACGCGTGTCTGACCTCGAGT
|||||||
NNTGCGCACAGCTGGAGCTCAGATCT (SEQ ID NO:6)

Five hundred pmol of adaptor were added to the tag DNA in a total volume of 50 µl of 1× ligase buffer containing 10U of T4 DNA ligase and incubated overnight at 16 °C. The ligase was inactivated by incubation at 65 °C for 10 min.

Step 5: PCR amplification of the tags

PCR amplification of the tags was carried out using sense and antisense primers designed to match the two adaptor sequences.

The following primers were used:

forward

5' TCTAGACTCGAGGTCGACACGC (SEQ ID NO:7)

and reverse

5' GCAGTGGTATCAACGCAGAGTCC (SEQ ID NO:8)

Step 6: Tag concatenation

The PCR product was electrophoresed on a polyacrylamide gel to isolate the 85 bp tag band. After phenol:CIA extraction and ethanol precipitation, the DNA was suspended in TE (pH 7.5). DNA was ligated with TA

U.S. PATENT APPLICATION OF SAMAL ET AL.

cloning vector (In Vitrogen, Inc, Carlsbad, CA).

Transformation was carried out according to the protocol provided by the manufacturer.

Transformed *E. coli* cells were grown in 100 ml of ampicillin-containing Terrific Broth at 37 °C, shaken at 300 rpm for 16 hr. Plasmid DNA preparation was carried out using Maxi kit (Qiagen Inc). About 750 µg DNA was obtained which was suspended in 500 µl of water.

The digestion of the purified plasmid DNA was carried out in a volume of 750 µl using 2 Units of Mlu I per µg of plasmid DNA for 4 hours. The resulting 28 bp tags were purified by electrophoresis on a 1.0% agarose gel in TAE buffer.

The 28 bp band was cut out of the gel, and eluted using a freeze-thaw technique. The DNA was extracted with phenol:CIA and ethanol precipitated in the presence of 4 µl glycogen and 100 µl of 10 M ammonium acetate per every 300 µl of sample. DNA was then resuspended in 16 µl water.

Concatemers were formed in a final volume of 20 µl using 1 µl of T4 DNA ligase (NEB, 400 units/µl). Concatemers were fractionated on an agarose gel isolating greater than 500 bp fragments. The fragments were purified using the Qiaex (Qiagen, Valencia, CA) protocol following

U.S. PATENT APPLICATION OF SAMAL ET AL.

the manufactures's instructions. The large molecular weight concatemers were then ligated into Mlu I-digested, alkaline phosphatase-treated, pBlueScript plasmid in which an Mlu I site had been engineered.

Results

The accuracy with which one can align a short cDNA sequence to the genomic sequence depends upon the length of the cDNA sequence. This is illustrated in TABLE 1 below. Using the NCBI Database of 47,584 known and hypothetical mRNAs, short expressed sequences (tags) from the 5' end of mRNAs were extracted and aligned against the genomic database. The result clearly demonstrates that at least 20 bases and preferably 32 bases or more of a contiguous sequence of mRNA are required to obtain a unique genomic match and thereby to identify a coding region from a genomic database.

U.S. PATENT APPLICATION OF SAMAL ET AL.

TABLE 1: Effect of Tag Length on Unique Genomic Hits

TAG LENGTH	%TAGS WITH UNIQUE GENOMIC HIT
14	5.76
16	37.56
18	74.47
20	84.56
32	89.44
36	90.07
40	90.61

However, currently, there is no enzyme, which can reproducibly generate 20 bases or longer fragments of double stranded cDNAs. We have developed a method to generate such expressed fragments. By obtaining one or more successive shorter fragments (tags) of 10-20 bases, which can then be aligned against the genomic sequence, the method generates two tandem tags which, in effect, produces a long contiguous sequence of 20 bases or greater. As a preferred embodiment, we have used an enzyme, Bpm I, which generates 16 base pair tags each time and 32 base pair tandem tags when aligned. A schematic outline of the method is shown in FIGURES 1A, 1B and 1C.

As an example, a tandem tag library, i.e., two successive tag libraries from a single cDNA sample, was generated from the mRNA of a human cortical neural stem

U.S. PATENT APPLICATION OF SAMAL ET AL.

cell culture consisting of approximately 2×10^7 cells. The resulting tag libraries were sequenced, aligned against the human genomic database, and pairs of tags, which align perfectly end to end on the genomic sequence were identified as tandem tags. Some of the tandem tags are shown in TABLE 2 and TABLE 3.

In TABLE 2, the two tandem 16-mer tags which uniquely and perfectly match known mRNA sequences are shown. The NCBI database of 47,584 known and hypothetical mRNAs was used as the template. In TABLE 3, the human genomic database was used first as the template to generate tandem tags. These were then compared to the mRNA database to verify whether the tandem tags indeed identified a coding region. These tandem tags are also found to be tandem within a known mRNA. BLAST of mRNA sequence to the human genome reveals that tandem genomic alignment was correct in each case.

TABLE 2: Examples of 16mer Tags Found to be Tandem within Known Transcripts

TAGS	MATCHING mRNA ACCESSION NO.	TANDEM TAG SEQUENCE POS.	mRNA NAME/DESCRIPTION
GCGCGGTGTGGTGGCA (SEQ ID NO:9) / GCAGGCGCAGCCCCAGC (SEQ ID NO:10)	NM_001024.2	14	Homo sapiens ribosomal protein S21 (RPS21), mRNA
GATAGATCGCCATCAT (SEQ ID NO:11) / GAACGACACCCGTAAC (SEQ ID NO:12)	NM_033022.1	24	Homo sapiens ribosomal protein S24 (RPS24), mRNA
TAGATCGCCATCATGA (SEQ ID NO:13) / ACGACACCCGTAAC (SEQ ID NO:14)	NM_033022.1	26	Homo sapiens ribosomal protein S24 (RPS24), mRNA
CTGCGGTGGAGCCGCC (SEQ ID NO:15) / ACCAAAATGCAGATT (SEQ ID NO:16)	NM_002954.2	23	Homo sapiens ribosomal protein S27a (RPS27A), mRNA
GTGGAGCTGTCGCCAT (SEQ ID NO:17) / GAAGGTCGAGCTGTGC (SEQ ID NO:18)	NM_000986.1	26	Homo sapiens ribosomal protein L24 (RPL24), mRNA
GCCATCGTGTGTGTT (SEQ ID NO:19) / CTTGACTCCGCTGCTC (SEQ ID NO:20)	NM_001000.1	3	Homo sapiens ribosomal protein L39 (RPL39), mRNA

CAGCACCATGCGGTT (SEQ ID NO:21) / GGCAAGAACAAGCGCC (SEQ ID NO:22)	NM_001006.1	30	Homo sapiens ribosomal protein S3A (RPS3A), mRNA
CTTGAACCTGGAGGC (SEQ ID NO:23) / GGAGGTTGCAGTGAAC (SEQ ID NO:24)	XM_040175.1	2779	Homo sapiens NADH dehydrogenase (ubiquinone) Fe-S protein 8 (23kD) (NADH-coenzyme Q reductase) (NDUFS8), mRNA
CTTGAACCCAGGAGGT (SEQ ID NO:25) / GGAGGTTGCAGTGATC (SEQ ID NO:26)	XM_035578.1	1853	Homo sapiens similar to X-like 1 protein (LOC91023), mRNA
GTGTGTGTGTGTGTGT (SEQ ID NO:27) / GTTTGTGTGTGTGTGT (SEQ ID NO:28)	NM_016352.1	2513	Homo sapiens carboxypeptidase A3 (LOC51200), mRNA

TABLE 3: Examples of Tags with Tandem Genome Alignment and Tandem mRNA Alignment; mRNA CDS found at location of Tandem Genome Alignment

TAGS	GENOME LOCATION OF TANDEM MATCH	MRNA LOCATION OF TANDEM MATCH	BLAST RESULTS OF MRNA TO GENOME ALIGNMENT
CTGCGGTGGAGCCGCC (SEQ ID NO:29) / ACCAAAATGCAGATT (SEQ ID NO:30)	NT_007741.6 MINUS strand @ 1,292,008	NM_002954.2 @ 23 (RPS27a)	NT_007741.6 MINUS strand 1,292,043 - 1,291,507
GTGGAGCTGTCGCCAT (SEQ ID NO:31) / GAAGTCGAGCTGTGC (SEQ ID NO:32)	NT_007592.6 MINUS strand @ 1,993,254	NM_000986.1 @ 26 (RPL24)	NT_007592.6 MINUS strand 1,993,292 - 1,992,861
GCCATCGTGTGTGTT (SEQ ID NO:33) / CTTGACTCCGCTGCTC (SEQ ID NO:34)	NT_007236.6 MINUS strand @ 3,673,626	NM_001000.1 @ 3 (RPL39)	NT_007236.6 MINUS strand 3,673,641 - 3,673,273
CAGCACCATGGCGGTT (SEQ ID NO:35) / GGCAAGAACAAGCGCC (SEQ ID NO:36)	NT_007816.6 MINUS strand @ 2,168,098 and 2,229,441	NM_001006.1 @ 30 (RPS3A)	NT_007816.6 MINUS strand 2,168,129 - 2,167,273 NT_007816.6 MINUS strand 2,229,472 - 2,228,616
CTTGAACCCAGGAGGT (SEQ ID NO:37) / GGAGTTGCAGTGATC (SEQ ID NO:38)	NT_010204.6 PLUS strand @ 1,527,899	XM_035578.1 @ 1853 (X-like 1 protein)	NT_010204.6 PLUS strand 1,472,273 - 1,527,982

NT_029281.1 PLUS

CTTGAACCCAGGAGGT (SEQ ID NO:39) /	NT_029281.1 PLUS strand @ 84,817	XM_043233.1 @ 875 (AK022192)	NT_029281.1 PLUS strand 83,943 - 86,079
TGCAGTGAGCCAAGAT (SEQ ID NO:40)			

U.S. PATENT APPLICATION OF SAMAL ET AL.

To further test the efficiency of the tandem tags to identify coding regions within the human genome, 400 random 16-mers from the first tag library and 400 random 16-mers from the second tag library were selected. Tandem tags were identified from the genomic database. As shown in TABLE 4, the 32-mer tandem tags were vastly more efficient in zeroing on the uniquely matching coding region of the human genome than the individual 16-mer tags.

U.S. PATENT APPLICATION OF SAMAL ET AL.

TABLE 4: Tandem vs. Non-tandem Efficiency

TAGS	GENOME MATCHES
GCACTTTGGGAGGCCGGCTCACGCCTGTAATC (SEQ ID NO:41)	1
GCACTTTGGGAGGCCG (SEQ ID NO:42)	157,201
GCTCACGCCTGTAATC (SEQ ID NO:43)	170,672
CACGCCCCGTAATCCCAAGCACTTTGGGAGGCT (SEQ ID NO:44)	1
CACGCCCCGTAATCCCA (SEQ ID NO:45)	1,337
AGCACTTTGGGAGGCT (SEQ ID NO:46)	132,561
AGCACTTTGGGAGGCTGAGATCGAGACCATCC (SEQ ID NO:47)	2
AGCACTTTGGGAGGCT (SEQ ID NO:48)	132,561
GAGATCGAGACCATCC (SEQ ID NO:49)	66,177
GCTTGAACCTGGGAGGGGAGGTTGCAGTGAGC (SEQ ID NO:50)	10
GCTTGAACCTGGGAGG (SEQ ID NO:51)	62,182
GGAGGTTGCAGTGAGC (SEQ ID NO:52)	162,173
GGCCAACATGGCGAAACCCGTCTCTACTAAAA (SEQ ID NO:53)	47
GGCCAACATGGCGAAA (SEQ ID NO:54)	17,111
CCCGTCTCTACTAAAA (SEQ ID NO:55)	138,143
GTGGAGCTTGCAGTGAGCCGAGATCGCGCCAC (SEQ ID NO:56)	1180
GTGGAGCTTGCAGTGA (SEQ ID NO:57)	14,992
GCCGAGATCGCGCCAC (SEQ ID NO:58)	20,598

U.S. PATENT APPLICATION OF SAMAL ET AL.

The key notion that two 16-mer tags can be aligned against the genomic database to identify a unique 32-mer coding sequence was further tested in silico in the following analysis. Using the set of 13,904 Unique RefSeq known mRNAs, two consecutive 16-mer tags were extracted near the 5' end of 1,000 mRNAs. These 16-mer tags were then pooled into a single "bin" to mimic a tag library. We then asked whether we could successfully recover, first, the tandem tags, and, second, the correct coding region by aligning the individual 16-mer tags against the human genome database. The 32 bp result set of tandem genome alignments was compared to the original 1,000 32bp known mRNA tandem. The results are summarized in TABLE 5 below.

Approximately 75% of the 32-mer sequences could be recovered by the tandem method. The remaining 25% not found in the genome are most likely due to the gaps and incomplete sequences present in the current version of the human genome database. The false positives, which appear because two 16-mer tags paired up illegitimately, constituted about 2%.

TABLE 5: In silico validation of the tandem tag method

TEST #	mRNA 32-MER SET	mRNA 16-MER SET	32-MER GENOME ALIGNMENTS	DISTINCT 32-MER TANDEMIS	32-MER mRNAs FOUND	GENOME FALSE POSITIVES
1	1000 (995 distinct)	2000 (1988 distinct)	35,874	727	720 (720/995 = 72.4%)	7 (7/727 = 0.96%)
2	1000 (991 distinct)	2000 (1982 distinct)	5,513	746	728 (728/991 = 73.5%)	18 (18/746 = 2.41%)
3	1000 (993 distinct)	2000 (1981 distinct)	154,854	758	752 (752/993 = 75.7%)	6 (6/758 = 0.79%)
4	1000 (992 distinct)	2000 (1981 distinct)	175,420	778	770 (770/992 = 77.6%)	8 (8/778 = 1.03%)
5	1000 (990 distinct)	2000 (1979 distinct)	910	736	729 (729/990 = 73.6%)	7 (7/736 = 0.95%)
6	1000 (992 distinct)	2000 (1984 distinct)	2,642	759	739 (739/992 = 74.5%)	20 (20/759 = 2.64%)
7	1000 (991 distinct)	2000 (1982 distinct)	1,436	735	730 (730/991 = 73.6%)	5 (5/735 = 0.68%)
8	1000 (992 distinct)	2000 (1983 distinct)	184,449	753	742 (742/992 = 74.8%)	11 (11/753 = 1.46%)
AVG 1-8	992 distinct sets	1983 distinct tags		749	74.5%	1.365%
9	3000 (2960 dist.)	6000 (5913 distinct)	177,607	2266	2212 (2212/2960 = 4.7%)	54 (54/2266 = 2.38%)

U.S. PATENT APPLICATION OF SAMAL ET AL.

Tags once extracted from the sequenced concatemers are usually subjected to a clustering protocol to positively match the tags to known transcripts or to the human genome. This is done due to the redundant occurrence of some of the 16 base pair tags within the genome, which does not allow the mining novel gene transcripts. Since the first set of tags and their tandem tags are generated from undefined ends of double-stranded cDNAs, each transcript is highly likely to generate multiple overlapping or closely spaced tags. Also, the number of such tags per transcript should be proportional to the relative abundance of the transcript in the sample. By aligning all tags against mRNA database and/or against the human genome, a stretch of physical sequence of the corresponding transcript is identified.

An example of a clustering protocol is shown below. Prior to clustering, 16 bp tags were extracted from sequenced concatemers and aligned to FASTA files of human genome, mRNA, and EST sequence databases. The output from this alignment program yields an alignment table for each respective sequence database. Each row in the alignment table is an exact location where one of the tags was found

U.S. PATENT APPLICATION OF SAMAL ET AL.

in the sequence database (GenBank Accession, strand, sequence position).

Using the genome or mRNA alignment table, tag hits are clustered by scanning each sequence (genome contig or mRNA) to group tags that are proximal to each other. The clustering program accepts two criteria: maximum hit-to-hit distance and minimum number of tag hits needed to define a cluster. The program picks up the first tag alignment and places it into the cluster bin. It continues down the genome strand until it finds the next alignment. If its distance away from the last alignment placed in the cluster bin is less than the maximum hit-to-hit distance then it is placed in the cluster bin. Clustering is finished when the next hit is too far away or the program finishes scanning the genome contig strand. If the number of hits in the cluster bin are at least the minimum number set by the user, then a cluster is created and the program outputs to a table the cluster location and other relevant information. With an mRNA alignment table, the cluster program works exactly the same way except that it scans down each mRNA instead of a genomic contig.

U.S. PATENT APPLICATION OF SAMAL ET AL.

To ensure high quality clusters, in this example, a maximum hit-to-hit distance of no greater than the tag length (hits must be adjacent or overlapping) was used. Minimum cluster size was 3 hits.

TAG CLUSTER EXAMPLES

1) Clustering against mRNA transcript database

(Refseq + Genome Annotation mRNAs)

CLUST ID	GENBGI	BEGIN POS	END POS	NUM TAGS
1	450185 8	1821	1846	6

mRNA ID:

>gi|4501858|ref|NM_001609.1| Homo sapiens acyl-Coenzyme A dehydrogenase, short/branched chain (ACADSB), nuclear gene encoding mitochondrial protein, mRNA
(2682 bp)

Location of transcript in Genome:

NT_008926.7|17472331 PLUS strand

64789 - 64929 (1003 - 1143)
66802 - 66906 (1142 - 1246)
*67437 - 68879 (1243 - 2682)

NT_027097.4 PLUS strand

1770323 - 1770376 (4 - 57)
1795662 - 1795822 (57 - 217)
1799051 - 1799154 (215 - 318)

U.S. PATENT APPLICATION OF SAMAL ET AL.

*matching genome cluster should be:

68015 - 68040 (1821 - 1846).

Clustering against Human Genome database:

CLUSTID	GENBGI	STRAND	BEGINPOS	ENDPOS	NUMTAGS
3411961	17472331	PLUS	68015	68040	6

This corresponds with expected cluster location and size.

2) mRNA Cluster(s):

CLUSTID	GENBGI	BEGINPOS	ENDPOS	NUMTAGS
2	4502010	1364	1396	8
3	4502010	1533	1562	7
4	4502010	1587	1623	8

>gi|4502010|ref|NM_000476.1| Homo sapiens adenylate kinase

1 (AK1), mRNA

(2271 bp)

mRNA matches Genome:

NT_029366.3|17449540 MINUS strand

1803682 - 1803643 (1 - 40)
1800671 - 1800631 (41 - 81)
1799083 - 1799043 (80 - 120)
1798874 - 1798709 (117 - 282)
1797960 - 1797843 (281 - 398)
1794533 - 1794339 (398 - 592)
*1794098 - 1792410 (589 - 2271)

*matching genome clusters should be:

1793291 - 1793323 (1396 - 1364)
1793125 - 1793150 (1562 - 1533)
1793064 - 1793100 (1623 - 1587)

U.S. PATENT APPLICATION OF SAMAL ET AL.

Genome Cluster(s):

CLUSTID	GENBGI	STRAND	BEGINPOS	ENDPOS	NUMTAGS
1862419	17449540	MINUS	1793062	1793098	8
1862420	17449540	MINUS	1793124	1793153	7
1862422	17449540	MINUS	1793289	1793321	8

3) mRNA Cluster(s):

CLUSTID	GENBGI	BEGINPOS	ENDPOS	NUMTAGS
5	4502042	1927	1959	9
6	4502042	2010	2047	6
7	4502042	2058	2131	12

>gi|4502042|ref|NM_000694.1| Homo sapiens aldehyde
dehydrogenase 3 family, member B1 (ALDH3B1), mRNA
(2790 bp)

mRNA matches Genome:

NT_008940.7|17472907 PLUS strand

1472982 - 1473028 (1 - 47)
1477929 - 1478094 (44 - 209)
1481160 - 1481272 (208 - 321)
1481406 - 1481528 (320 - 442)
1481798 - 1481889 (436 - 527)
1482346 - 1482431 (525 - 610)
1484116 - 1484504 (607 - 996)
1485227 - 1485398 (996 - 1167)
1488638 - 1488743 (1160 - 1265)
*1490381 - 1491906 (1263 - 2790)

*matching genome cluster(s) should be:

1491045 - 1491077 (1927 - 1959)
1491128 - 1491165 (2010 - 2047)
1491176 - 1491249 (2058 - 2131)

U.S. PATENT APPLICATION OF SAMAL ET AL.

Genome Cluster(s):

CLUSTID	GENBGI	STRAND	BEGINPOS	ENDPOS	NUMTAGS
3473301	17472907	PLUS	1491044	1491076	9
3473302	17472907	PLUS	1491127	1491164	6
3473303	17472907	PLUS	1491175	1491248	12

4) mRNA Cluster(s):

CLUSTID	GENBGI	BEGINPOS	ENDPOS	NUMTAGS
7012	14786455	2347	2385	12

>gi|14786455|ref|XM_009672.4| Homo sapiens

phosphoenolpyruvate carboxykinase 1 (soluble) (PCK1), mRNA

(2642 letters)

mRNA matches Genome:

NT_011362.7|17484369 PLUS strand
 21189036 - 21189118 (1 - 83)
 21189283 - 21189548 (80 - 345)
 21189983 - 21190167 (345 - 529)
 21190607 - 21190812 (526 - 731)
 21190941 - 21191128 (732 - 919)
 21191447 - 21191642 (919 - 1084)
 21192080 - 21192307 (1081 - 1308)
 21192394 - 21192529 (1307 - 1442)
 21192952 - 21193049 (1439 - 1536)
 21193261 - 21194369 (1534 - 2642)

*matching genome cluster(s) should be:

21194074 - 21194112 (2347 - 2385)

Genome Cluster(s):

CLUSTID	GENBGI	STRAND	BEGINPOS	ENDPOS	NUMTAGS
4332399	17484369	PLUS	21194074	21194112	12

U.S. PATENT APPLICATION OF SAMAL ET AL.

5) mRNA Cluster(s):

CLUSTID	GENBGI	BEGINPOS	ENDPOS	NUMTAGS
647	5174710	1385	1410	5
648	5174710	1446	1484	10

>gi|5174710|ref|NM_005992.1| Homo sapiens T-box 1 (TBX1),

transcript variant B, mRNA

(1538 bp)

mRNA matches Genome:

NT_011519.9|17484914 PLUS strand

2892106 - 2892148 (1 - 43)
 2894958 - 2895080 (41 - 163)
 2896306 - 2896684 (162 - 540)
 2898641 - 2898747 (537 - 643)
 2899557 - 2899729 (641 - 813)
 2900361 - 2900516 (814 - 969)
 2901160 - 2901229 (969 - 1038)
 2901304 - 2901406 (1037 - 1139)
 2918314 - 2918438 (1137 - 1261)
 *2918714 - 2918996 (1256 - 1538)

*matching genome cluster(s) should be:

2918843 - 2918868 (1385 - 1410)
 2918904 - 2918942 (1446 - 1484)

Genome Cluster(s):

CLUSTID	GENBGI	STRAND	BEGINPOS	ENDPOS	NUMTAGS
4343636	17484914	PLUS	2918843	2918868	5
4343637	17484914	PLUS	2918904	2918942	10

Occasionally, alignment of two tandem 16-mer tags on the human genome produced false 32-mer sequences that

U.S. PATENT APPLICATION OF SAMAL ET AL.

probably do not exist in real transcripts. These represent a false-pairing against the human genome and are false-positives. Such false pairing can be reduced by using a second 5' adaptor containing two degenerate nucleotide bases. This example is shown below:

Bpm I digestion

5' ... C T G G A G (N)16^ ... 3'
3' ... G A C C T C (N)14^ ... 5' (SEQ ID NO:59)

The first adaptor:

GCAGTGGTATCAACGCAGAGTCCACGCGTCTGGAG
||||||||||||||||||||||||||||||
CACCATAGTTGCGTCTCAGGTGCGCAGACCTC_p (SEQ ID NO:3)

The second adaptor with 2 nn on the 3' end of the first strand:

GCAGTGGTATCAACGCAGAGTCCACGCGTCTGGAGNN
||||||||||||||||||||||||||||||
CACCATAGTTGCGTCTCAGGTGCGCAGACCTC_p (SEQ ID NO:60)

Bpm I digestion leaves 3'-overhang of two nucleotides on the bottom strands of the leftover cDNA to which the second adaptor with two nn 3' overhang on the top strand is ligated. These two nucleotides are conserved in the second tag after second Bpm I cut. Hence the last two nucleotides of the first tag and the first two nucleotides of the

20090205-030603

U.S. PATENT APPLICATION OF SAMAL ET AL.

'putative' tandem tag are the same. This prevents the random matching of all the available tags to the first tag and decreases significantly the artificial combination between two random 16 mers.

TABLE 6 below lists other type II restriction enzymes that generate short DNA fragments away from the recognition sites and could be used in this method.

U.S. PATENT APPLICATION OF SAMAL ET AL.

TABLE 6: Type II restriction enzymes with asymmetric recognition sequences:

Type II restriction enzymes

Cuts after 4n	Ear I, Sap I,
Cuts after 5n	Alw I, Bmr I, Bsa I, BsmA I, BsmB I, MlyI, PleI,
Cuts after 6n	Bbs I, BciV I, Fau I,
Cuts after 7n	Mnl I,
Cuts after 8n	Aar I, BfuA I, BspM I, Hph I, Mbo II, SspD5I, Sth132I,
Cuts after 9n	SfaN I,
Cuts after 10n	BseR I, BspCN I, Hga I,
Cuts after 11n	AceIII, Eci I, TaqII, Tth111II,
Cuts after 12n	Bbv I, RleAI,
Cuts after 13n	BcefI, Fok I
Cuts after 14n	BceA I, BsmF I, StsI,
Cuts after 16n	Bce83I, Bpm I, Bsg I, Eco57I, Eco57MI,
Cuts after 20n	MmeI

U.S. PATENT APPLICATION OF SAMAL ET AL.

While the invention has been described in connection with what is presently considered to be the most practical and preferred embodiments, it is to be understood that the invention is not limited to the disclosed embodiments, but on the contrary is intended to cover various modifications and equivalent arrangements included within the spirit and scope of the appended claims.

Thus, it is to be understood that variations in the present invention can be made without departing from the novel aspects of this invention as defined in the claims. All patents and articles cited herein are hereby incorporated by reference in their entirety and relied upon.